

**Peringkasan Multidokumen Otomatis dengan Menggunakan
Log-Likelihood Ratio (LLR) dan *Maximal Marginal
Relevance* (MMR) untuk Artikel dengan Topik Penyakit
Menular Bahasa Indonesia**

SKRIPSI

**Diajukan untuk Memenuhi Salah Satu Syarat Mencapai Gelar Strata
Satu Program Studi Informatika**



Disusun oleh :

**IKHWAN NIZWAR AKHMAD
M0511025**

**PROGRAM STUDI INFORMATIKA
FAKULTAS MATEMATIKA DAN ILMU PENGETAHUAN ALAM
UNIVERSITAS SEBELAS MARET
SURAKARTA
2017**

SKRIPSI

Peringkasan Multidokumen Otomatis dengan *Menggunakan Log-Likelihood Ratio (LLR)* dan *Maximal Marginal Relevance (MMR)* untuk Artikel dengan Topik Penyakit Menular Bahasa Indonesia

Disusun Oleh :

IKHWAN NIZWAR AKHMAD

M0511025

Skripsi ini telah disetujui untuk dipertahankan di hadapan dewan penguji pada tanggal,

Pembimbing I

Pembimbing II

**Drs. Bambang Harjito M.App.Sc., Ph.D
NIP. 19621130 199103 1 002**

**Dr. Eng. Anto Satriyo Nugroho
NIP. 19701021 198911 1 001**

SKRIPSI

Peringkasan Multidokumen Otomatis dengan Menggunakan Log-Likelihood Ratio (LLR) dan Maximal Marginal Relevance (MMR) untuk Artikel dengan Topik Penyakit Menular Bahasa Indonesia

Disusun Oleh:

IKHWAN NIZWAR AKHMAD

M0511025

**Telah dipertahankan dihadapan Dewan Penguji
pada tanggal, 15 Desember 2017**

Susunan Dewan Penguji

- 1. Drs. Bambang Harjito M.App.Sc., Ph.D. (Ketua) ()
NIP. 196211301991031002**
- 2. Dr. Anto Satriyo Nugroho, M. Eng. (Sekretaris) ()
NIP. 197010211989111001**
- 3. Sari Widya Sihwi, S.Kom., M.T.I. (Anggota) ()
NIP. 198304122009122003**
- 4. Denis Eka Cahyani, S.Kom, M.Kom. (Anggota) ()
NIP. 1991031020161001**

Disahkan Oleh

Kepala Program Studi Informatika

**Drs. Bambang Harjito M.App.Sc., Ph.D.
NIP. 196211301991031002**

HALAMAN MOTTO

“Bila engkau ingin satu, maka jangan ambil dua. Karena satu menggenapkan, tapi dua melenyapkan.” - Dee

HALAMAN PERSEMBAHAN

Untuk Abi, Umik, Hani, Iin, dan Sekar.

KATA PENGANTAR

Puji syukur kehadiran Allah SWT yang telah melimpahkan rahmat-Nya sehingga penulis dapat menyelesaikan skripsi dengan judul “Peringkasan Multidokumen Otomatis dengan Menggunakan *Log-Likelihood Ratio* (LLR) dan *Maximal Marginal Relevance* (MMR) untuk Artikel dengan Topik Penyakit Menular Bahasa Indonesia”. Skripsi ini disusun sebagai salah satu syarat dalam memperoleh gelar Sarjana Komputer pada Program Studi Informatika Universitas Sebelas Maret.

Dalam proses penelitian dan penyusunan skripsi ini, penulis telah mendapatkan banyak bantuan dan dukungan dari berbagai pihak. Untuk itu, penulis menyampaikan terima kasih kepada:

1. Bapak, Ibu, serta Adik yang telah dan terus-menerus memberikan dukungan material dan non-material selama proses penyelesaian tugas akhir.

2. Bapak Dr. Anto Satriyo Nugroho, M.Eng selaku Pembimbing Tugas Akhir yang telah memberikan banyak solusi atas permasalahan penulis selama proses penyelesaian tugas akhir, serta membuka wawasan penulis mengenai ranah penelitian di Indonesia.

3. Bapak Drs. Bambang Harjito, M.APP.Sc, Ph.D. selaku Kepala Program Studi Informatika Universitas Sebelas Maret, Pembimbing Akademik, dan Pembimbing Tugas Akhir, yang telah memberikan banyak motivasi, saran, dan bimbingan selama proses studi dan penyelesaian tugas akhir.

4. Bapak Prof. Ir. Ari Handono Ramelan, M.Sc.(Hons), Ph.D. selaku Dekan Fakultas MIPA Universitas Sebelas Maret.

5. Seluruh Dosen Program Studi Informatika yang telah membagikan begitu banyak dasar ilmu pengetahuan dan pengalamannya kepada penulis.

6. Teman-teman Informatika dan Ngarsapura Creative Media, yang telah menjadi bagian dari perjalanan studi, pembentukan karakter diri, juga susah senang sebagai mahasiswa dan rekan kerja.

Kesempurnaan memang hanya milik Allah SWT. Namun Penulis juga menyadari bahwa penelitian dan skripsi ini masih jauh dari penelitian dan laporan yang ideal. Meskipun demikian, penulis sangat berharap skripsi ini bisa bermanfaat bagi pembaca.

Surakarta, November 2017

Penulis

Peringkasan Multidokumen Otomatis dengan Menggunakan *Log-Likelihood Ratio* (LLR) dan *Maximal Marginal Relevance* (MMR) untuk Artikel dengan Topik Penyakit Menular Bahasa Indonesia

Ikhwan Nizwar Akhmad

Program Studi Informatika, Fakultas Matematika dan Ilmu Pengetahuan
Alam, Universitas Sebelas Maret

ABSTRAK

Peningkatan jumlah informasi yang tersedia di internet disamping memberikan manfaat, juga memunculkan masalah tersendiri. Mesin pencarian modern sudah cukup baik untuk mendapatkan informasi tertentu. Namun jumlah informasi yang sangat banyak terkadang menyebabkan pencari informasi kesulitan mendapatkan intisari dari informasi yang dicari. Kondisi yang disebabkan oleh terlalu banyaknya informasi yang tersedia dikenal sebagai *information overload*.

Peringkasan multidokumen otomatis merupakan salah satu solusi untuk masalah ini. Meskipun metode peringkasan multidokumen otomatis sudah dikembangkan sejak 20 tahun lalu, penerapannya dalam Bahasa Indonesia masih terbatas. Artikel dengan topik penyakit menular merupakan salah satu studi kasus yang menarik untuk peringkasan multidokumen Bahasa Indonesia. Informasi mengenai penyakit menular dibutuhkan oleh masyarakat sehingga terdapat banyak informasi mengenai topik ini di internet. Kondisi ini menyebabkan kemungkinan *information overload* untuk pencarian dalam topik ini.

Dalam penelitian ini, akan diterapkan metode peringkasan multidokumen otomatis dengan menggunakan *Log-Likelihood Ratio* (LLR) untuk mendapatkan *topic signature*, dan *Maximal Marginal Relevance* pada artikel dengan topik penyakit menular untuk mendapatkan ringkasan dengan sedikit perulangan informasi.

Penelitian ini menghasilkan ringkasan dengan nilai akurasi sebesar 0,4 (dengan menggunakan ROUGE-S9). Selain itu, dalam penelitian ini didapatkan bahwa *topic signature* (beserta akurasinya) memegang peran penting dalam proses peringkasan dokumen otomatis.

Kata Kunci: peringkasan multidokumen otomatis, *topic signature generation*

Automatic Multidocument Summarization utilizing Log-Likelihood Ratio (LLR) and Maximal Marginal Relevance (MMR) for Infectious Diseases Articles in Indonesian

Ikhwan Nizwar Akhmad

*Department of Informatics, Mathematics and Natural Science Faculty,
Sebelas Maret University*

ABSTRACT

Increasing number of information that available on the Internet, along with its benefit, also comes with its own problems. Modern search engines are smart enough to bring the most corresponding information, but the immense number of information provided sometimes bring more confusion than clarity. This condition is known as information overload.

Automatic multidocument summarization is a way to overcome this particular problem. But despite being heavily researched more than 20 years ago, its implementations for Bahasa Indonesia are limited. Articles about infectious disease is one of the ideal case study for multidocument summarization for Bahasa Indonesia. Information about infectious disease are essential for general public therefore many information about it is available on the Internet. This condition could trigger information overload when someone do an internet search in this topic.

In this research, we try to implement multidocument summarization technique for articles with infectious disease topic in Bahasa Indonesia utilizing Log Likelihood Ratio (LLR) to obtain topic signatures and Maximal Marginal Relevance (MMR) to generate relevant summary with minimal information redundancy.

Our summarization method generated a summary with 0.4 fmeasure using ROUGE-S9 evaluation. Also, we found that topic signature (with its accuracy) takes an important role on generating good summaries.

Keywords: *multidocument summarization, topic signatures generation*

DAFTAR ISI

Halaman Motto.....	iv
Halaman Persembahan.....	v
Kata Pengantar.....	vi
Abstrak.....	viii
<i>Abstract</i>	ix
Daftar Isi.....	x
Daftar Gambar.....	xii
Daftar Tabel.....	xiii
BAB 1 Pendahuluan.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	4
1.3 Batasan Masalah.....	4
1.4 Tujuan Penelitian.....	4
1.5 Manfaat Penelitian.....	5
1.6 Sistematika Penulisan.....	5
BAB 2 Tinjauan Pustaka.....	6
2.1 Dasar Teori.....	6
2.1.1 Peringkasan Dokumen Otomatis.....	6
2.1.2 <i>Text Preprocessing</i>	9
2.1.3 <i>Log Likelihood Ratio (LLR)</i>	12
2.1.4 <i>Maximal Marginal Relevance (MMR)</i>	14
2.1.5 Evaluasi Metode Peringkasan Dokumen Otomatis.....	16
2.2 Penelitian Terkait.....	19
BAB 3 Metodologi Penelitian.....	26
3.1 Studi Literatur.....	26
3.2 Pengumpulan Data.....	27
3.3 Metode yang Diusulkan.....	27
3.3.1 <i>Text Preprocessing</i>	28
3.3.2 <i>Topic Signature Generation</i> dengan LLR.....	29
3.3.3 Peringkasan Multidokumen Otomatis dengan MMR.....	29

3.4 Metode Evaluasi.....	29
BAB 4 Hasil dan Pembahasan.....	32
4.1 <i>Topic Signature Generation</i> dengan LLR.....	32
4.2 Pemilihan Kalimat Pertama Ringkasan.....	35
4.3 Pemilihan Kalimat Ringkasan dengan menggunakan MMR.....	37
4.4 Evaluasi Metode Peringkasan.....	42
BAB 5 Penutup.....	44
5.1 Kesimpulan.....	44
5.2 Saran.....	45
Daftar Pustaka.....	46

DAFTAR GAMBAR

Gambar 2.1 Contoh <i>case folding</i> pada kalimat	10
Gambar 2.2 Contoh <i>tokenization</i> pada paragraf	11
Gambar 2.3 Contoh <i>stopword removal</i>	12
Gambar 2.4 Contoh <i>stemming</i> pada kumpulan kata	12
Gambar 2.5 Metode evaluasi peringkasan otomatis	17
Gambar 3.1 Metodologi Penelitian	26
Gambar 3.2 Metode peringkasan multidokumen otomatis yang diusulkan. 28	
Gambar 3.3 Evaluasi metode peringkasan yang diusulkan.....	30
Gambar 4.1 <i>Histogram</i> nilai LLR untuk setiap kata dalam korpus.....	34
Gambar 4.2 ROUGE-S9 <i>F-Measure</i>	43

DAFTAR TABEL

Tabel 2.1 <i>Contingency Table</i>	13
Tabel 2.2 Penelitian Terkait.....	22
Tabel 4.1 <i>Contingency table</i> untuk kata "infeksi"	31
Tabel 4.2 Kata-kata dengan LLR tertinggi.....	34
Tabel 4.3. Pemilihan Kalimat Pertama Ringkasan.....	35
Tabel 4.4. Contoh hasil ringkasan otomatis.....	39
Tabel 4.5. Ringkasan yang telah Diurutkan.....	41